

## CS 2750 Machine Learning Lecture 18

### Ensemble methods:

- **Mixture of experts**
- **Bagging**

Milos Hauskrecht

[milos@cs.pitt.edu](mailto:milos@cs.pitt.edu)

5329 Sennott Square

---

CS 2750 Machine Learning

## Ensemble methods

**Train and use multiple base models for either classification or regression problems**

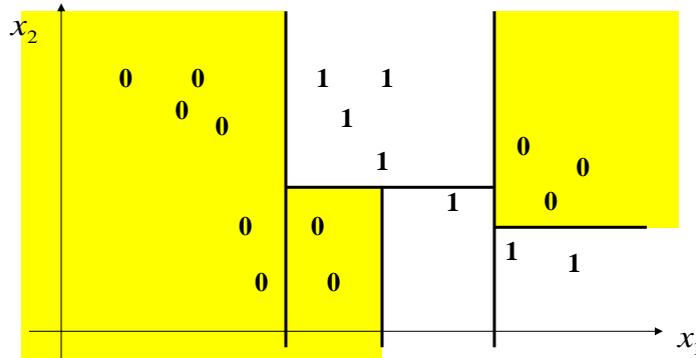
- **Mixture of experts**
  - each ‘base’ models (classifier, regressor) covers a different part (region) of the input space
  - All models are trained together on the same training data
- **Committee machines:**
  - each ‘base’ model (classifier, regressor) covers the complete input space
  - Each base model is trained on a slightly different train set
  - Combine predictions of all models to produce the output
    - **Goal:** Improve the accuracy of the ‘base’ model
  - **Methods: Bagging, Boosting, Stacking** (not covered)

---

CS 2750 Machine Learning

## Reviewing Decision trees

- An approach to classification that:
  - **Partitions the input space to regions**
  - **Classifies independently in every region**



CS 2750 Machine Learning

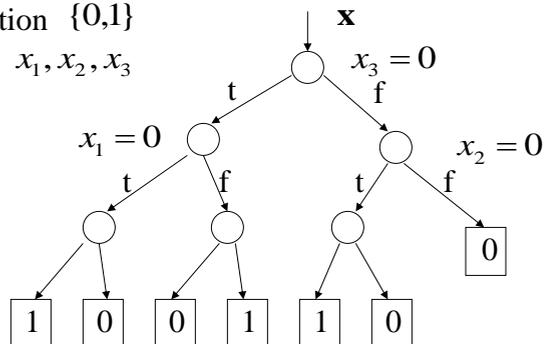
## Decision trees

- The partitioning idea is used in the **decision tree model**:
  - Split the space recursively according to inputs in  $\mathbf{x}$
  - Classify (assign class label) at the bottom of the tree

### Example:

Binary classification  $\{0,1\}$

Binary attributes  $x_1, x_2, x_3$



CS 2750 Machine Learning

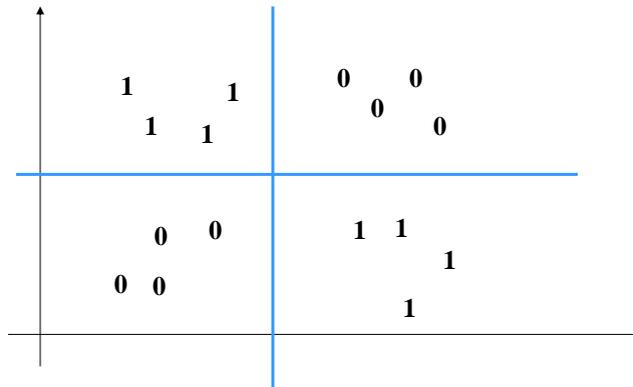
## Decision tree learning

- **Greedy learning algorithm:**
  - Repeat until no or small improvement in the purity
    - Find the attribute with the highest gain
    - Add the attribute to the tree and split the set accordingly
- Builds the tree in the top-down fashion
  - Gradually expands the leaves of the partially built tree
- The method is greedy
  - It looks at a single attribute and gain in each step
  - May fail when the combination of attributes is needed to improve the purity (parity functions)

CS 2750 Machine Learning

## Limitations of Decision trees

- **Greedy learning methods:** a combination of two or more attributes improves the impurity
- **Rectangular regions**

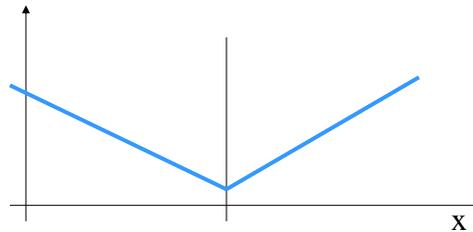


CS 2750 Machine Learning

## Mixture of experts model

- **Ensamble methods:**
  - Use a combination of simpler learners/model to improve their predictions
- **Mixture of expert model:**
  - Different input regions covered with different learners
  - A “soft” switching between learners

- **Mixture of experts**  
**Expert = learner**

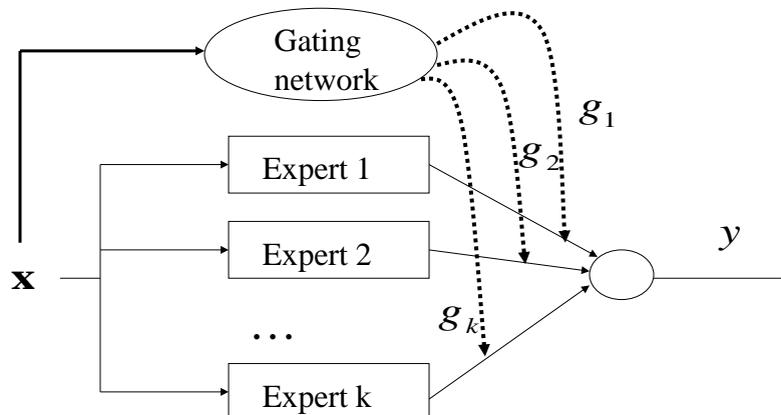


CS 2750 Machine Learning

## Mixture of experts model

- **Gating network** : decides what expert to use

$g_1, g_2, \dots, g_k$  - gating functions



CS 2750 Machine Learning

## Learning mixture of experts

- **Learning consists of two tasks:**
  - Learn the parameters of individual expert networks
  - Learn the parameters of the gating (switching) network
    - Decides where to make a split
- **Assume:** gating functions give probabilities
 
$$0 \leq g_1(\mathbf{x}), g_2(\mathbf{x}), \dots, g_k(\mathbf{x}) \leq 1 \quad \sum_{u=1}^k g_u(\mathbf{x}) = 1$$
- Based on the probability we partition the space
  - partitions belongs to different experts
- How to model the gating network?
  - **A multi-way classifier model:**
    - **softmax model**

CS 2750 Machine Learning

## Learning mixture of experts

- Assume we have a **set of linear experts**

$$\mu_i = \boldsymbol{\theta}_i^T \mathbf{x} \quad (\text{Note: bias terms are hidden in } \mathbf{x})$$
- Assume a **softmax gating network**

$$g_i(\mathbf{x}) = \frac{\exp(\boldsymbol{\eta}_i^T \mathbf{x})}{\sum_{u=1}^k \exp(\boldsymbol{\eta}_u^T \mathbf{x})} \approx p(\omega_i | \mathbf{x}, \boldsymbol{\eta})$$
- Likelihood of  $y$  (linear regression – assume errors for different experts are normally distributed with the same variance)

$$\begin{aligned}
 P(y | \mathbf{x}, \Theta, \boldsymbol{\eta}) &= \sum_{i=1}^k P(\omega_i | \mathbf{x}, \boldsymbol{\eta}) p(y | \mathbf{x}, \omega_i, \Theta) \\
 &= \sum_{i=1}^k \left[ \frac{\exp(\boldsymbol{\eta}_i^T \mathbf{x})}{\sum_{j=1}^k \exp(\boldsymbol{\eta}_j^T \mathbf{x})} \right] \left[ \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\|y - \mu_i\|^2}{2\sigma^2}\right) \right]
 \end{aligned}$$

CS 2750 Machine Learning

## Learning mixture of experts

### Learning of parameters of expert models:

#### On-line update rule for parameters $\theta_i$ of expert $i$

- If we know the expert that is responsible for  $\mathbf{x}$

$$\theta_{ij} \leftarrow \theta_{ij} + \alpha_{ij} (y - \mu_i) x_j$$

- If we do not know the expert

$$\theta_{ij} \leftarrow \theta_{ij} + \alpha_{ij} h_i (y - \mu_i) x_j$$

$h_i$  - responsibility of the  $i$ th expert = a kind of posterior

$$h_i(\mathbf{x}, y) = \frac{g_i(\mathbf{x}) p(y | \mathbf{x}, \omega_i, \theta)}{\sum_{u=1}^k g_u(\mathbf{x}) p(y | \mathbf{x}, \omega_u, \theta)} = \frac{g_i(\mathbf{x}) \exp(-1/2 \|y - \mu_i\|^2)}{\sum_{u=1}^k g_u(\mathbf{x}) \exp(-1/2 \|y - \mu_u\|^2)}$$

$g_i(\mathbf{x})$  - a prior       $\exp(\dots)$  - a likelihood

CS 2750 Machine Learning

## Learning mixtures of experts

### Learning of parameters of the gating/switching network:

- On-line learning of gating network parameters  $\eta_i$

$$\eta_{ij} \leftarrow \eta_{ij} + \beta_{ij} (h_i(\mathbf{x}, y) - g_i(\mathbf{x})) x_j$$

- The learning with conditional mixtures can be extended to learning of parameters of an **arbitrary expert network**
  - e.g. logistic regression, multilayer neural network

$$\theta_{ij} \leftarrow \theta_{ij} + \beta_{ij} \frac{\partial l}{\partial \theta_{ij}}$$
$$\frac{\partial l}{\partial \theta_{ij}} = \frac{\partial l}{\partial \mu_i} \frac{\partial \mu_i}{\partial \theta_{ij}} = h_i \frac{\partial \mu_i}{\partial \theta_{ij}}$$

CS 2750 Machine Learning

## Learning mixture of experts

**EM algorithm** for learning the mixture components

**Algorithm:**

Initialize parameters  $\Theta$

Repeat

Set  $\Theta' = \Theta$

1. **Expectation step**

$$Q(\Theta | \Theta') = E_{H|\mathbf{X},\mathbf{Y},\Theta'} \log P(\mathbf{H}, \mathbf{Y} | \mathbf{X}, \Theta, \xi)$$

2. **Maximization step**

$$\Theta = \arg \max_{\Theta} Q(\Theta | \Theta')$$

until no or small improvement in  $Q(\Theta | \Theta')$

- **Hidden variables are identities of expert networks responsible for (x,y) data points**

CS 2750 Machine Learning

## EM for Learning mixture of experts

- Assume we have a **set of linear experts**

$$\mu_i = \boldsymbol{\theta}_i^T \mathbf{x}$$

- Assume a **softmax gating network**

$$g_i(\mathbf{x}) = P(\omega_i | \mathbf{x}, \boldsymbol{\eta})$$

- **Q function to optimize**

$$Q(\Theta | \Theta') = E_{H|\mathbf{X},\mathbf{Y},\Theta'} \log P(\mathbf{H}, \mathbf{Y} | \mathbf{X}, \Theta, \xi)$$

- **Assume:**

–  $l$  indexes different data points

–  $\delta_i^l$  an indicator variable for the data point  $l$  to be covered by an expert  $i$

$$Q(\Theta | \Theta') = \sum_l \sum_i E(\delta_i^l | \mathbf{x}^l, y^l, \Theta', \boldsymbol{\eta}) \log(P(y^l, \omega_i | \mathbf{x}^l, \Theta, \boldsymbol{\eta}))$$

CS 2750 Machine Learning

## Learning mixture of experts with EM

- Assume:

- $l$  indexes different data points
- $\delta_i^l$  an indicator variable for data point  $l$  and expert  $i$

$$Q(\Theta | \Theta') = \sum_l \sum_i E(\delta_i^l | \mathbf{x}^l, y^l, \Theta', \boldsymbol{\eta}') \log(P(y^l, \omega_i | \mathbf{x}^l, \Theta, \boldsymbol{\eta}))$$

### Responsibility of the expert $i$ for $(x, y)$

$$E(\delta_i^l | \mathbf{x}^l, y^l, \Theta', \boldsymbol{\eta}') = h_i^l(\mathbf{x}^l, y^l) = \frac{g_i(\mathbf{x}^l) p(y | \mathbf{x}^l, \omega_i, \boldsymbol{\theta}')}{\sum_{u=1}^k g_u(\mathbf{x}^l) p(y^l | \mathbf{x}^l, \omega_u, \boldsymbol{\theta}')}$$

$$Q(\Theta | \Theta') = \sum_l \sum_i h_i^l(\mathbf{x}^l, y^l) \log(P(y^l, \omega_i | \mathbf{x}^l, \Theta, \boldsymbol{\eta}))$$

CS 2750 Machine Learning

## EM for learning the mixture of experts

- The maximization step boils down to the problem that is equivalent to the problem of finding the ML estimates of the parameters of the expert and gating networks

$$Q(\Theta | \Theta') = \sum_l \sum_i h_i^l(\mathbf{x}^l, y^l) \log(P(y^l, \omega_i | \mathbf{x}^l, \Theta, \boldsymbol{\eta}))$$

$$\log(P(y^l, \omega_i | \mathbf{x}^l, \Theta, \boldsymbol{\eta})) = \log P(y^l | \omega_i, \mathbf{x}^l, \Theta) + \log P(\omega_i | \mathbf{x}^l, \boldsymbol{\eta})$$

Expert network  $i$   
(Linear regression)

Gating network  
(Softmax)

- Note that any optimization technique can be applied in this step

CS 2750 Machine Learning

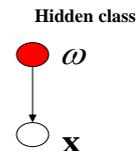
## Learning mixture of experts

- Note that we can use different expert and gating models
- For example:
  - Experts: logistic regression models

$$y_i = 1 / (1 + \exp(-\boldsymbol{\theta}_i^T \mathbf{x}))$$

- Gating network: a **generative latent variable model**

$$g_i(\mathbf{x}) = P(\omega_i | \mathbf{x}, \boldsymbol{\eta})$$



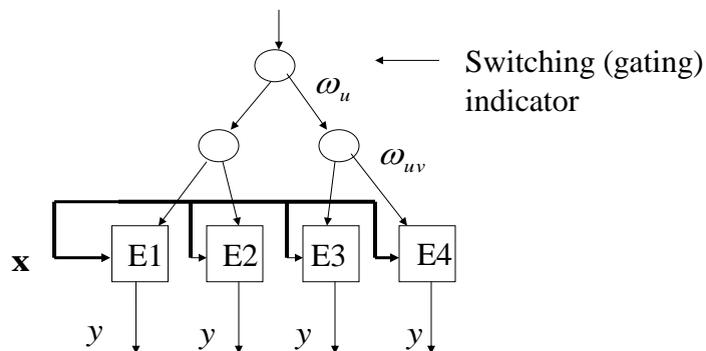
- Likelihood of  $y$ :

$$P(y | \mathbf{x}, \Theta, \boldsymbol{\eta}) = \sum_{u=1}^k P(\omega_u | \mathbf{x}, \boldsymbol{\eta}) p(y | \mathbf{x}, \omega_u, \Theta)$$

CS 2750 Machine Learning

## Hierarchical mixture of experts

- **Mixture of experts**: define a probabilistic split
- The idea can be extended to a **hierarchy of experts** (a kind of a probabilistic decision tree)



CS 2750 Machine Learning

## Hierarchical mixture model

An output is conditioned (gated) on multiple mixture levels

$$P(y | \mathbf{x}, \Theta) = \sum_u P(\omega_u | \mathbf{x}, \eta) \sum_v p(\omega_{uv} | \mathbf{x}, \omega_u, \xi_u) \dots \sum_s P(\omega_{uv..s} | \mathbf{x}, \omega_u, \omega_{uv}, \dots) \underline{P(y | \mathbf{x}, \omega_u, \omega_{uv}, \dots, \theta_{uv..s})}$$

**Individual experts**

- **Define**  $\Omega_{uv..s} = \{\omega_u, \omega_{uv}, \dots, \omega_{uv..s}\}$

$$P(\Omega_{uv..s} | \mathbf{x}, \Theta) = P(\omega_u | \mathbf{x}) P(\omega_{uv} | \mathbf{x}, \omega_u) \dots P(\omega_{uv..s} | \mathbf{x}, \omega_u, \omega_{uv}, \dots)$$

- **Then**

$$P(y | \mathbf{x}, \Theta) = \sum_u \sum_v \dots \sum_s P(\Omega_{uv..s} | \mathbf{x}, \Theta) P(y | \mathbf{x}, \Omega_{uv..s}, \Theta)$$

- **Mixture model is a kind of soft decision tree model**  
- **with a fixed tree structure !!**

CS 2750 Machine Learning

## Hierarchical mixture of experts

- Multiple levels of probabilistic gating functions

$$g_u(\mathbf{x}) = P(\omega_u | \mathbf{x}, \Theta) \qquad g_{v|u}(\mathbf{x}) = P(\omega_{uv} | \mathbf{x}, \omega_u, \Theta)$$

- Multiple levels of responsibilities

$$h_u(\mathbf{x}, y) = P(\omega_u | \mathbf{x}, y, \Theta) \qquad h_{v|u}(\mathbf{x}, y) = P(\omega_{uv} | \mathbf{x}, y, \omega_u, \Theta)$$

- How they are related?

$$\begin{aligned}
 \text{responsibility} \\
 P(\omega_{uv} | \mathbf{x}, y, \omega_u, \Theta) &= \frac{P(y | \mathbf{x}, \omega_u, \omega_{uv}, \Theta) P(\omega_{uv} | \mathbf{x}, \omega_u, \Theta)}{\sum_v P(y | \mathbf{x}, \omega_u, \omega_{uv}, \Theta) P(\omega_{uv} | \mathbf{x}, \omega_u, \Theta)} \\
 &= \sum_v P(y, \omega_{uv} | \mathbf{x}, \omega_u, \Theta) = P(y | \mathbf{x}, \omega_u, \Theta)
 \end{aligned}$$

CS 2750 Machine Learning

## Hierarchical mixture of experts

- **Responsibility for the top layer**

$$h_u(\mathbf{x}, y) = P(\omega_u | \mathbf{x}, y, \Theta) = \frac{P(y | \mathbf{x}, \omega_u, \Theta)P(\omega_u | \mathbf{x}, \Theta)}{\sum_u P(y | \mathbf{x}, \omega_u, \Theta)P(\omega_u | \mathbf{x}, \Theta)}$$

- But  $P(y | \mathbf{x}, \omega_u, \Theta)$  is computed while computing  $h_{v|u}(\mathbf{x}, y) = P(\omega_{uv} | \mathbf{x}, y, \omega_u, \Theta)$

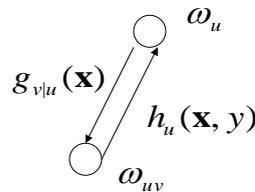
- **General algorithm:**

- **Downward sweep; calculate**

$$g_{v|u}(x) = P(\omega_{uv} | \mathbf{x}, \omega_u, \Theta)$$

- **Upward sweep; calculate**

$$h_u(\mathbf{x}, y) = P(\omega_u | \mathbf{x}, y, \Theta)$$



CS 2750 Machine Learning

## On-line learning

- Assume linear experts  $\mu_{uv} = \boldsymbol{\theta}_{uv}^T \mathbf{x}$

- **Gradients (vector form):**

$$\frac{\partial l}{\partial \boldsymbol{\theta}_{uv}} = h_u h_{v|u} (y - \mu_{uv}) \mathbf{x}$$

$$\frac{\partial l}{\partial \boldsymbol{\eta}} = (h_u - g_u) \mathbf{x} \quad \text{Top level (root) node}$$

$$\frac{\partial l}{\partial \boldsymbol{\xi}} = h_u (h_{v|u} - g_{v|u}) \mathbf{x} \quad \text{Second level node}$$

- Again: can it can be extended to different expert networks

CS 2750 Machine Learning

## Ensemble methods

- **Mixture of experts**
  - Multiple ‘base’ models (classifiers, regressors), each covers a different part (region) of the input space
- **Committee machines:**
  - Multiple ‘base’ models (classifiers, regressors), each covers the complete input space
  - Each base model is trained on a slightly different train set
  - Combine predictions of all models to produce the output
    - **Goal:** Improve the accuracy of the ‘base’ model
  - **Methods:**
    - **Bagging**
    - **Boosting**
    - Stacking (not covered)

---

CS 2750 Machine Learning

## Bagging (Bootstrap Aggregating)

- **Given:**
  - Training set of  $N$  examples
  - A class of learning models (e.g. decision trees, neural networks, ...)
- **Method:**
  - Train multiple ( $k$ ) models on different samples (data splits) and average their predictions
  - Predict (test) by averaging the results of  $k$  models
- **Goal:**
  - Improve the accuracy of one model by using its multiple copies
  - Average of misclassification errors on different data splits gives a better estimate of the predictive ability of a learning method

---

CS 2750 Machine Learning

## Bagging algorithm

- **Training**

- In each iteration  $t, t=1, \dots, T$ 
  - Randomly sample with replacement  $N$  samples from the training set
  - Train a chosen “base model” (e.g. neural network, decision tree) on the samples

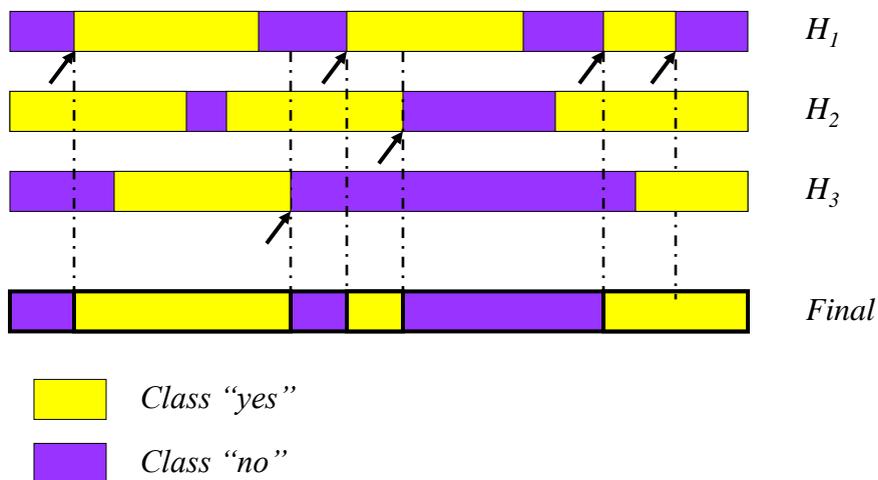
- **Test**

- For each test example
  - Start all trained base models
  - Predict by combining results of all  $T$  trained models:
    - **Regression:** averaging
    - **Classification:** a majority vote

CS 2750 Machine Learning

## Simple Majority Voting

Test examples



CS 2750 Machine Learning

## Analysis of Bagging

- **Expected error= Bias+Variance**

- *Expected error* is the expected discrepancy between the estimated and true function

$$E\left[\left(\hat{f}(X) - E[f(X)]\right)^2\right]$$

- *Bias* is squared discrepancy between *averaged* estimated and true function

$$\left(E[\hat{f}(X)] - E[f(X)]\right)^2$$

- *Variance* is expected divergence of the estimated function vs. its average value

$$E\left[\left(\hat{f}(X) - E[\hat{f}(X)]\right)^2\right]$$

CS 2750 Machine Learning

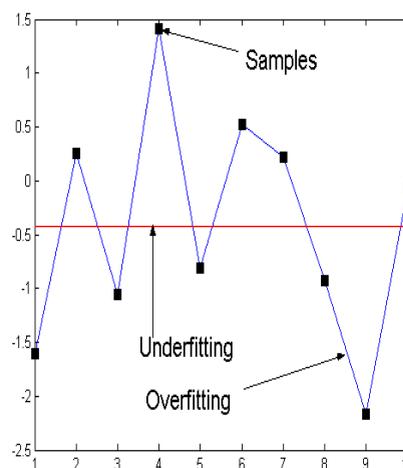
## When Bagging works? Under-fitting and over-fitting

- **Under-fitting:**

- High bias (models are not accurate)
- Small variance (smaller influence of examples in the training set)

- **Over-fitting:**

- Small bias (models flexible enough to fit well to training data)
- Large variance (models depend very much on the training set)



CS 2750 Machine Learning

## Averaging decreases variance

- **Example**
  - Assume we measure a random variable  $x$  with a  $N(\mu, \sigma^2)$  distribution
  - If only one measurement  $x_1$  is done,
    - The expected mean of the measurement is  $\mu$
    - Variance is  $\text{Var}(x_1) = \sigma^2$
  - If random variable  $x$  is measured  $K$  times ( $x_1, x_2, \dots, x_k$ ) and the value is estimated as:  $(x_1 + x_2 + \dots + x_k) / K$ ,
    - Mean of the estimate is still  $\mu$
    - But, variance is smaller:
      - $[\text{Var}(x_1) + \dots + \text{Var}(x_k)] / K^2 = K\sigma^2 / K^2 = \sigma^2 / K$
- Observe: **Bagging is a kind of averaging!**

CS 2750 Machine Learning

## When Bagging works

- **Main property of Bagging** (proof omitted)
  - Bagging **decreases variance** of the base model without changing the bias!!!
  - Why? averaging!
- **Bagging typically helps**
  - When applied with an **over-fitted base model**
    - High dependency on actual training data
- **It does not help much**
  - High bias. When the base model is robust to the changes in the training data (due to sampling)

CS 2750 Machine Learning

## Boosting

- **Mixture of experts**
  - One expert per region
  - Expert switching
- **Bagging**
  - Multiple models on the complete space, a learner is not biased to any region
  - Learners **are learned independently**
- **Boosting**
  - Every learner covers the complete space
  - Learners are biased to regions not predicted well by other learners
  - **Learners are dependent**